



Generative AI for Health Economic Evaluation

IHEA | Robert Smith, PhD | October 2024



rsmith@darkpeakanalytics.com



<https://github.com/dark-peak-analytics>



<https://www.linkedin.com/company/dark-peak-analytics>

Limited peer reviewed literature

arXiv:2407.11054 [cs.LG] Sept 21, 2024 Version 3

Generative AI for Health Technology Assessment: Opportunities, Challenges, and Policy Considerations

Rachael L. Fleurence, PhD, MSc¹, Jiang Bian, PhD^{1,2,3,4}, Xinyan Wang, PhD^{5,6}, Hsa Xu, PhD⁷, Dalis Dawoud, PhD^{8,9}, Mitch Higashi, PhD¹⁰, Jagpreet Chhatwal, PhD¹¹

¹Office of the Director, National Institute of Biomedical Imaging and Bioengineering, National Institutes of Health, Bethesda, MD, United States
²Health Outcomes and Biomedical Informatics, College of Medicine, University of Florida, Gainesville, FL, United States
³Health Outcomes and Biomedical Informatics, Clinical and Translational Science Institute, University of Florida, Gainesville, FL, United States
⁴Office of Data Science and Research Implementation, University of Florida Health, Gainesville, FL, United States
⁵Tulane University School of Public Health and Tropical Medicine, New Orleans, LA, United States
⁶Intelligent Medical Objects, Rosemont, IL, United States
⁷Department of Biomedical Informatics and Data Science, School of Medicine, Yale University, New Haven, CT, United States
⁸National Institute for Health and Care Excellence, London, United Kingdom
⁹Cairo University, Faculty of Pharmacy, Cairo, Egypt
¹⁰ISPOR - The Professional Society for Health Economics and Outcomes Research, Lawrenceville, NJ, United States
¹¹Institute for Technology Assessment, Massachusetts General Hospital, Harvard Medical School, Boston, MA, United States
¹²Center for Health Decision Science, Harvard University, Boston, MA, United States

Funding: Dr Dalis Dawoud reports partial funding from the European Union's Horizon 2020 research and innovation programme under Grant Agreement No. 82516 (Next Generation Health Technology Assessment (HTA)) project. No other funding was received.

Acknowledgements: The authors thank Dr Tala Fakheri for her comments on earlier versions of this manuscript.

1

Fleurence, R., Bian, J., Wang, X., Xu, H., Dawoud, D., Fakhouri, T., Higashi, M. and Chhatwal, J., 2024. Generative AI for Health Technology Assessment: Opportunities, Challenges, and Policy Considerations. arXiv preprint arXiv:2407.11054.

Pharmacoeconomics - Open (2024) 9:191–203
<https://doi.org/10.1007/s40264-024-00674-9>

ORIGINAL RESEARCH ARTICLE

Artificial Intelligence to Automate Health Economic Modelling: A Case Study to Evaluate the Potential Application of Large Language Models

Tim Reason¹, William Rawlinson¹, Julia Langham¹, Andy Gimblett¹, Bill Malcolm¹, Sven Klijn¹

Accepted: 1 February 2024 / Published online: 10 February 2024
 © The Author(s) 2024

Abstract
Background: Current generation large language models (LLMs) such as Generative Pre-Trained Transformer 4 (GPT-4) have achieved human-level performance on many tasks including the generation of computer code based on natural input. This study aimed to assess whether GPT-4 could be used to automatically programme two published health economic analyses.
Methods: The two analyses were partitioned survival models evaluating interventions in non-small cell lung cancer (NSCLC) and renal cell carcinoma (RCC). We developed prompts which instructed GPT-4 to programme the NSCLC and RCC models in R, and which provided descriptions of each model's methods, assumptions and parameter values. The results of the generated scripts were compared to the published values from the original, human-programmed models. The models were replicated 15 times to capture variability in GPT-4's output.
Results: GPT-4 fully replicated the NSCLC model with high accuracy: 100% (15/15) of the artificial intelligence (AI)-generated NSCLC models were error-free or contained a single minor error, and 93% (14/15) were completely error-free. GPT-4 closely replicated the RCC model, although human intervention was required to simplify an element of the model design (one of the model's fifteen input calculations because it used too many sequential steps to be implemented in a single prompt). With this simplification, 87% (13/15) of the AI-generated RCC models were error-free or contained a single minor error, and 40% (6/15) were completely error-free. Error-free model scripts replicated the published incremental cost-effectiveness ratios to within 1%.
Conclusion: This study provides a promising indication that GPT-4 can have practical applications in the automation of health economic model construction. Potential benefits include accelerated model development timelines and reduced costs of development. Further research is necessary to explore the generalisability of LLM-based automation across a larger sample of models.

1 Introduction

We are living through a golden age of innovations and the development of new treatments for many diseases. However, this is occurring at a time of increasing demand, primarily due to an ageing population with complex health needs, together with constrained healthcare resources and budgets. Health economic models, which provide evidence of the relative costs and benefits of new health technologies compared with existing technologies [1], are vital tools for informing health decision making, particularly health technology assessments that inform national decisions for market access and reimbursement [2].

Key Points for Decision Makers
 GPT-4, a current generation large language model (LLM), automatically replicated two published health economic models with high accuracy, based on instructions about how the models should be designed and what input values should be used.
 This is a promising early indication that LLMs could be used to automate building health economic models, which could reduce the costs of health economic analysis, accelerate model development timelines and reduce the risk of error in modelling.

© The Author(s)
tim.reason@nhs.uk
¹ Evidera Scientific, Middlesex, PO Box 10, London W12 7PP, UK
² Bristol Myers Squibb, Cheshire, UK
³ Bristol Myers Squibb, Princeton, NJ, USA

△ Adis

Reason, T., Rawlinson, W., Langham, J., Gimblett, A., Malcolm, B. and Klijn, S., 2024. Artificial Intelligence to Automate Health Economic Modelling: A Case Study to Evaluate the Potential Application of Large Language Models. *Pharmacoeconomics-Open*, 8(2), pp.191-203.

Wellcome Open Research

assertHE: an R package to improve quality assurance of HTA models

Robert A Smith^{1,2,3}, Yevgeny Samyushkin¹, Wael Mohammed^{1,2}, Felicity Lamrock⁴, Tom Ward^{1,5}, Jack Smith¹, Alan Martin¹, Paul Schneider^{1,6}, Dawn Lee¹, Gianluca Balio¹, Howard Thom¹, Nathan Green¹, Marina Richardson¹, Mohamed El Ailli^{1,2,3}, Xavier Pouwels¹, Calum Lewis^{1,2}, and Baris Deniz^{1,2,10}

¹Dark Peak Analytics, Sheffield, UK
²University of Sheffield, Sheffield, UK
³USK, London, UK
⁴University of Exeter Medical School, Exeter, EX1 2LU, UK
⁵PeritAG, University of Exeter Medical School, University of Exeter, St Luke's Campus, Exeter, EX1 2LU, UK
⁶Department of Biostatistics, University College London, London, UK
⁷Department of Population Health Sciences, Bristol Medical School, University of Bristol, Bristol, England, UK
⁸Section of Health Technology and Services Research, Technical School of Health, Faculty of Behavioural, Management, and Social Sciences, University of Twente, Enschede, Overijssel, The Netherlands
⁹National Science Research Centre, Queen's University Belfast, UK
¹⁰Current affiliation: Alia Solutions LLC, 2010 N Church St PEB 20087, Wilmington, Delaware 19802-4447, US
¹¹Institute for Clinical and Economic Review (ICER), Boston, MA, USA
¹²Department of Health Sciences, Faculty of Science, Vrije Universiteit Amsterdam, Amsterdam Public Health research institute, Amsterdam, The Netherlands
¹³National Health Care Institute (ZIN), Cluj-Napoca, The Netherlands

Corresponding author: Robert A. Smith robert@darkpeakanalytics.com
robert@darkpeakanalytics.com, HTA, Health Economics, Unit Testing, Model Validation, Model Development

Background: Health economic models are increasingly used to inform decisions about the allocation of healthcare resources. Ensuring the robustness and validity of these models is critical. Currently, quality assurance is conducted by both technical and non-technical experts assessing different components of the model manually. This is resource-intensive. Understanding how the different components of the model fit together is time-consuming, and testing every part of the model is sometimes not feasible in the time available. To aid in this, we have developed the assertHE R package.
Methods: The open source assertHE package provides testing functionality for three building and reviewing health economic models built in R programming language. It provides a series of checks which can be integrated into the model development workflow to reduce the probability of common errors. It also provides a suite of functions which allow users to better understand the network of functions contained in a model, where they are defined, if (and when) they are tested, and provides a simple means to quantify the extent to which they are tested.
Results: We applied the assertHE package to three open source health economic models. By allowing long lists of checks to be executed within the model code and how to visualise the network of functions, we saw test coverage, and about a 50% reduction in the number of errors. We used the assertHE package to generate a summary of any function in the code base. We have worked with collaborators from industry, regulators and academics to develop the package to be applicable to the widest possible range of models, making adaptations to the source code based upon feedback.
Conclusions: assertHE offers an open-source toolset for health economists building and reviewing models, promoting collaborative development and facilitating a more robust and efficient quality assurance process.

Page 1 of 10

https://www.courses.darkpeakanalytics.com/products/digital_downloads/asserthe-quality-assurance

Background & Motivation

Sure, I'll
review your
model

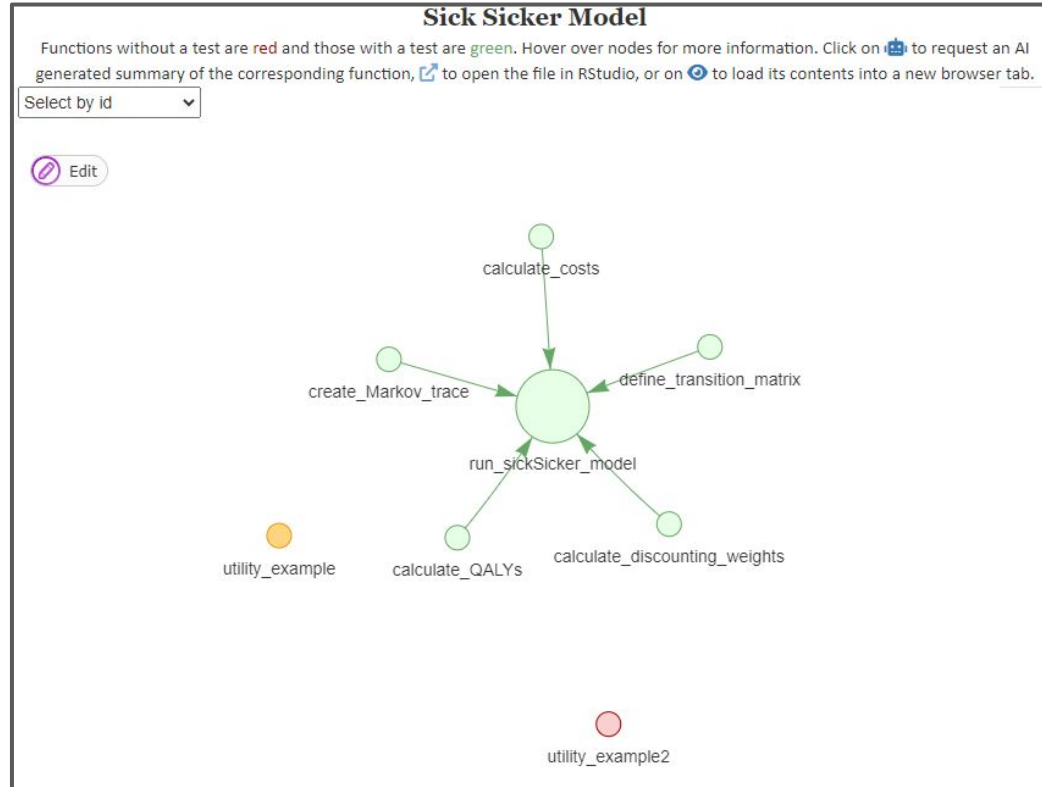
Good
Luck!

```
def PRINT_HEAD  
printf(" (%c%c%c)", c, ccharU(...func...));  
* ends;  
// D:\Scene(C, scene)V骨序pare_.clear();  
// (MoMoY骨序pareG.empty())  
#ifdef PC(C, o); // PTMoY骨序Pg(MoMoM+pare)(V+pareG[o]); // return;  
#else G=1; G=0;  
if(o->type == OB_ARMATURE)  
// V骨序V=LIB.resize(20);  
bArmature* arm骨=(bArmature*)o->data; int i=0;  
// Object*o_骨=o->c("o_骨", C); oG=o_骨; Object*oV=o->c("oV", C); Object*o骨=o->c("o骨", C);  
// o骨2=o->c("o骨2", C);  
bPoseChannel* pcA=CTX_data_active_pose_bone(C), *pcX=NULL, *pc骨=NULL;  
if(pcA&&pcA->parent&&pcA->parent->parent)  
Object*oTarget=o->c("oTarget", C); 正Vector丁目标=loc+骨(oTarget->ubmat);  
if(三("TypeOfIK", "PoseBone", pcA)!=1)return;  
骨序pare_骨=MoMoY骨序pareG[o];
```







assertHE

assertHE model reviewer



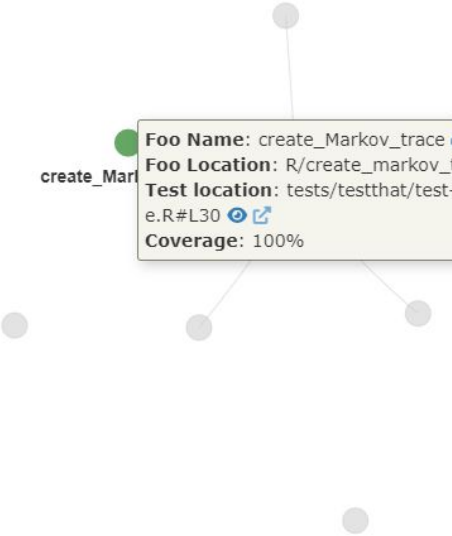
assertHE model reviewer


Sick Sicker Model



Functions without a test are **red** and those with a test are **green**. Hover over nodes for more information. Click on  to request an AI generated summary of the corresponding function,  to open the file in RStudio, or on  to load its contents into a new browser tab.



create_Markov_trace ▼

Edit



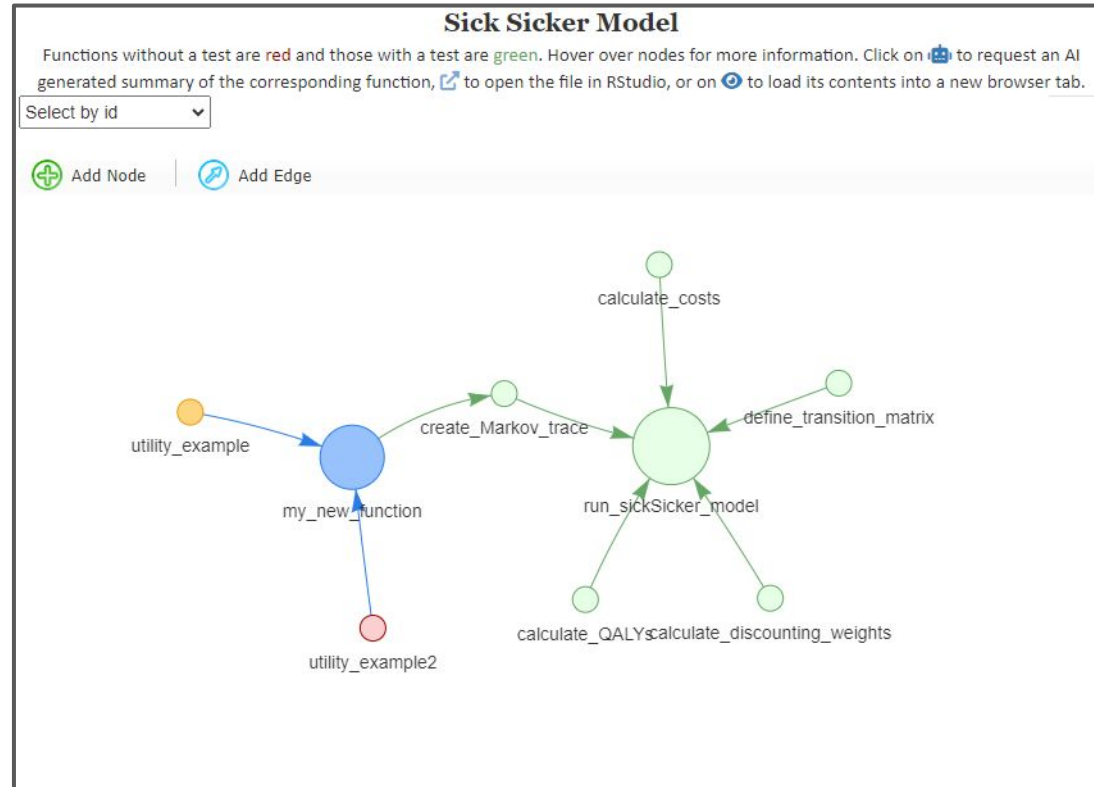
Foo Name: create_Markov_trace 

Foo Location: R/create_markov_trace.R#L45  

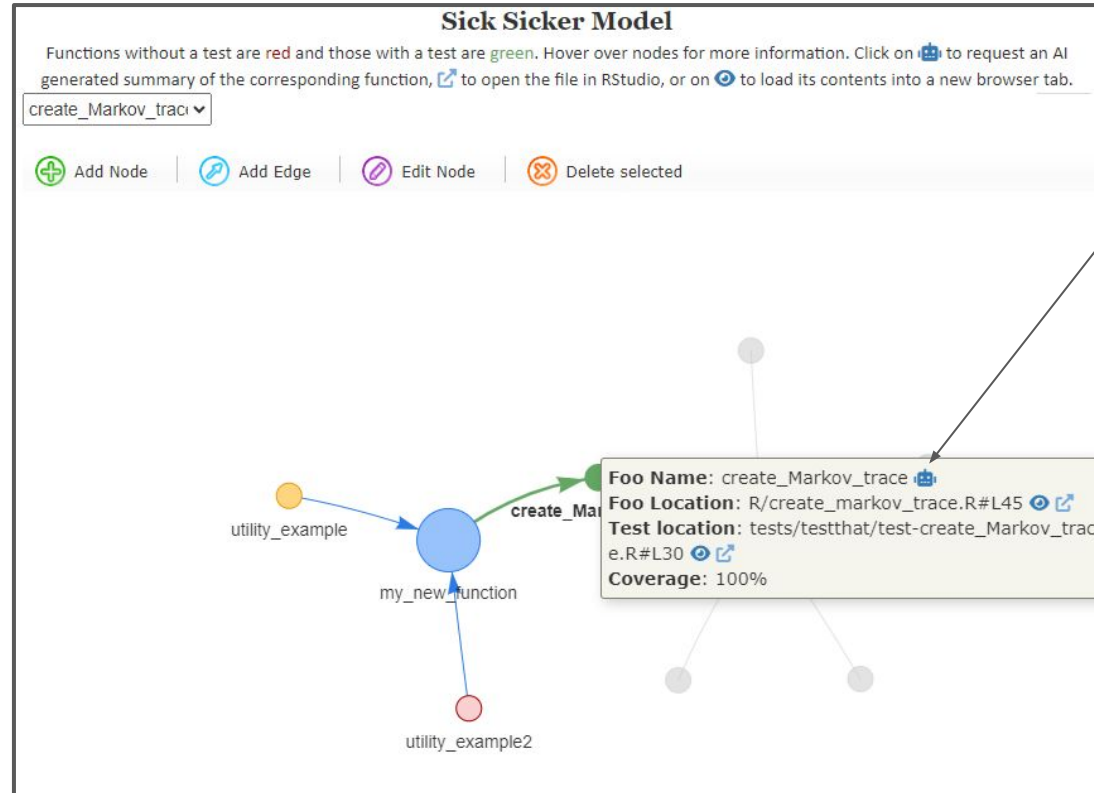
Test location: tests/testthat/test-create_Markov_trace.R#L30  

Coverage: 100%

assertHE model reviewer

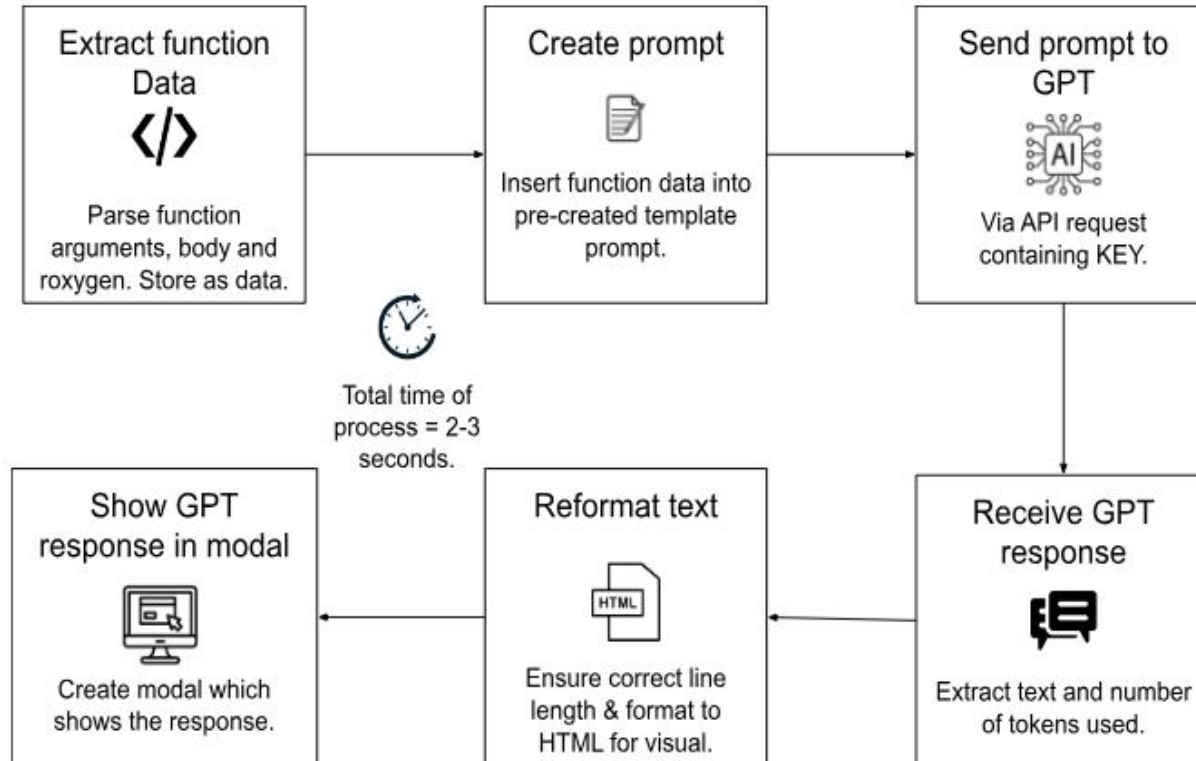


assertHE model reviewer

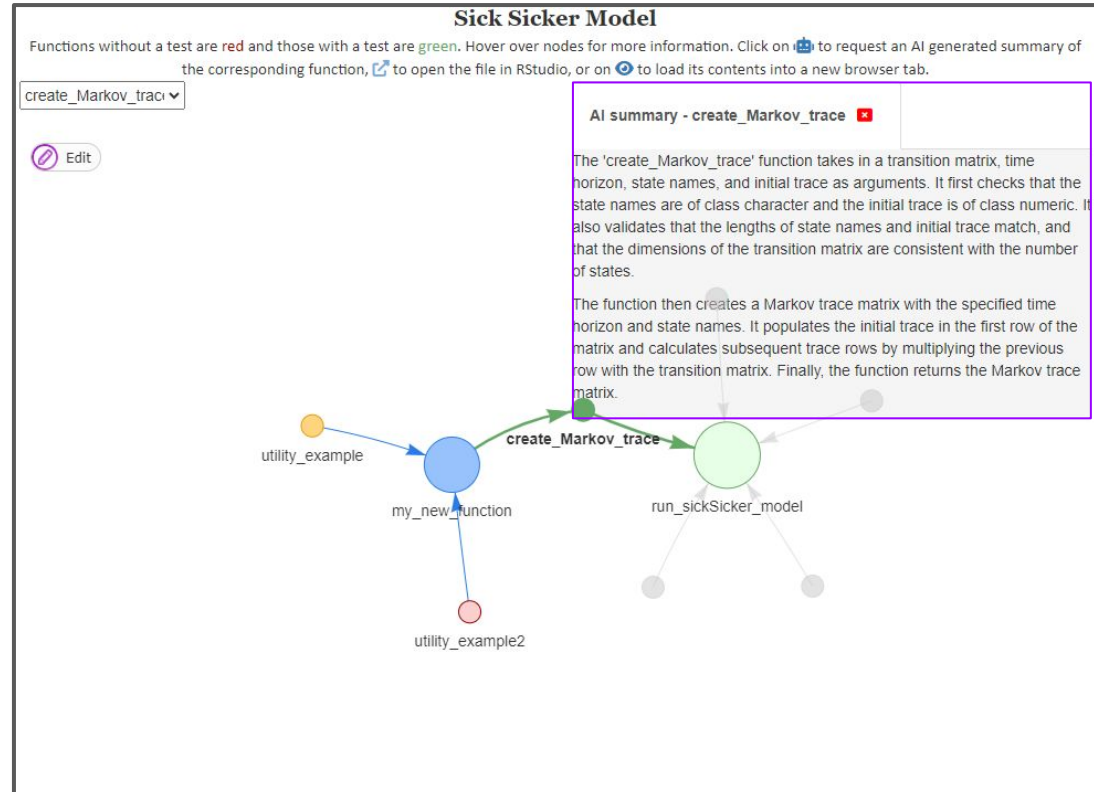


Generate LLM
summary of
function.

assertHE model reviewer



assertHE model reviewer








Case Studies



Case Study

Function Network

Functions without a test are **red** and those with a test are **green**. Hover over nodes for more information. Click on  to request an AI generated summary of the corresponding function,  to open the file in RStudio, or on  to load its contents into a new browser tab.

run_probsa ▼

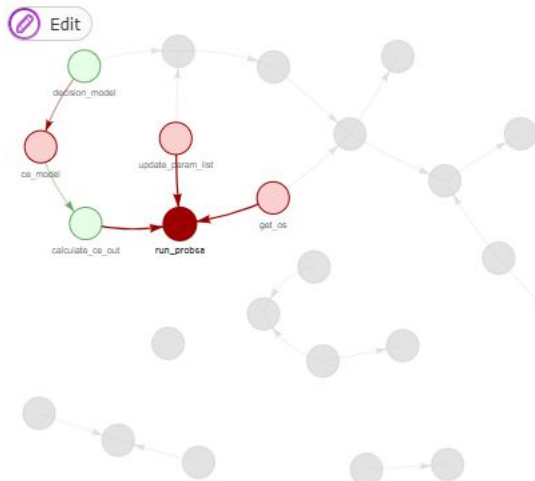
AI summary - run probsa ✕

```
'''html
```

The `'run_probsa'` function runs probabilistic sensitivity analysis (PSA) on a given input dataset. If the `'parallel'` argument is set to `TRUE`, the function parallelizes the PSA process using multiple cores based on the operating system. It then calculates costs and effects for each simulation, aggregates the results, and returns them in separate data frames.

If the 'parallel' argument is set to FALSE, the function runs the PSA simulations in series. It iterates through each simulation, updates parameters, calculates costs and effects, and prints the progress. Finally, it aggregates the results and returns them in separate data frames. The function returns a list containing the costs and effects data frames.

523


ELSEVIER

ScienceDirect
Contents lists available at www.sciencedirect.com
journal homepage: www.elsevier.com/locate/jl

Economic Evaluation

CDX2 Biomarker Testing and Adjuvant Therapy for Stage II Colon Cancer: An Exploratory Cost-Effectiveness Analysis

Fernando Alarid-Escufo, PhD, Deborah Scheig, MD, MPH, Karen M. Kuntz, ScD

ABSTRACT

Objective: Adjuvant chemotherapy is not recommended for patients with average-risk stage II (T3N0) colon cancer. Nevertheless, a subgroup of these patients who are CDX2-negative might benefit from adjuvant chemotherapy. We evaluated the cost-effectiveness of testing for the absence of CDX2 expression followed by adjuvant chemotherapy (fluorouracil combined with oxaliplatin [FOLFOX]) for patients with stage II colon cancer.

Methods: We developed a decision model to simulate a hypothetical cohort of 65-year-old patients with average-risk stage I colon cancer with 724 of these patients being CDX2-negative under 2 different interventions: (1) test for the absence of CDX2 expression followed by *adjuvant* chemotherapy for CDX2-negative patients and (2) no CDX2 testing and no *adjuvant* chemotherapy for any patient. We derived disease progression parameters, *adjuvant* chemotherapy effectiveness and utilities from published analyses, and cancer care costs from the Surveillance, Epidemiology, and End Results (SEER) Medicare data. Sensitivity analyses were conducted.

Results: Testing for CDX2 followed by FOLFIR for CDX2-negative patients had an incremental cost-effectiveness ratio of \$5500/quality-adjusted life-years (QALYs) compared with no CDX2 testing and no FOLFIR (6.874 vs 6.838 discounted QALYs and \$80,991 vs \$80,797 discounted US dollar lifetime costs). In sensitivity analyses, considering a cost-effectiveness threshold of \$100,000/QALY, testing for CDX2 followed by FOLFIR on CDX2-negative patients remains cost-effective for hazard ratios of <0.975 of the effectiveness of FOLFIR in CDX2-negative patients in reducing the rate of developing a metastatic recurrence.

Conclusions: Testing tumors of patients with stage II colon cancer for CD132 and administration of adjuvant treatment to the subgroup found CD132-negative is a cost-effective and high-value management strategy across a broad range of plausible assumptions.

Keywords: CD02, cost-effectiveness analysis, decision-analytic model, immunohistochemistry testing, stage II colon cancer

VALUE HEALTH. 2022; 25(3):409-419

Introduction

Adjuvant chemotherapy is not recommended for patients with average-risk stage II (T3N0) colon cancer,¹⁻⁴ and thus, these patients are usually treated with surgery alone.⁵ Nevertheless, a recent study by Dalerba et al⁶ described a small subgroup of patients with average-risk stage II colon cancer who lack expression of the CD32 transcription factor that associated with clinical benefit from adjuvant chemotherapy. CD32 is a master transcription factor involved in intestinal development⁷ and serves as a candidate biomarker of mature colon epithelial tissues. In this study, we used a novel, sensitive, and specific method to conduct a systematic search for a biomarker to identify undifferentiated tumors in a collection of human colon gene expression array experiments from the National Center for Biotechnology Information

National Cancer Institute, which were analyzed for CD22 expression by immunohistochemical (IHC) analysis. Among all tumors analyzed in the validation data set, 48 of 669 (7.2%) were CD22-negative, defined as completely lacking CD22 expression or showing expression in a minority of malignant epithelial cells.⁴⁴ The study also showed that CD22-negative patients had poorer 5-year disease-free survival (DFS) than CD22-positive patients (those with biomarker expression). More importantly, the 5-year DFS was greater for the CD22-negative patients who received adjuvant chemotherapy than similar patients who did not receive adjuvant chemotherapy. The ability to test average-risk patients versus high-risk patients for CD22 expression may allow for adjuvant chemotherapy to a subgroup most likely to benefit and reduce cancer mortality and minimize adjuvant chemotherapy harms.⁴⁵

This study aims to quantify the long-term benefits, costs, and cost-effectiveness of testing average-risk patients with stage II colon cancer for the absence of CTCE2 biomarker expression following

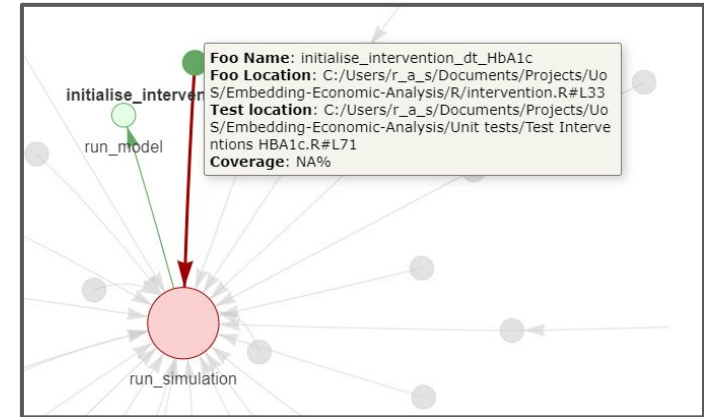
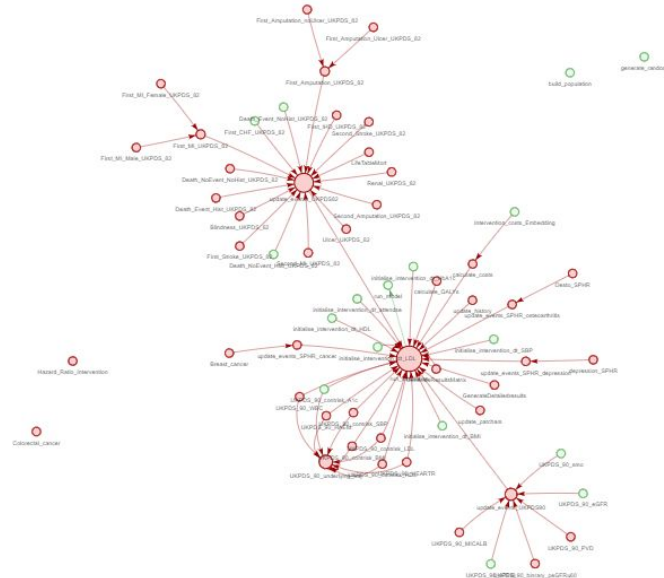
Case Study: Embedding Economic Analysis

Embedding-Economic-Analysis Repository

Functions without a test are **red** and those with a test are **green**. Hover over nodes for more information.

Select by id ▼

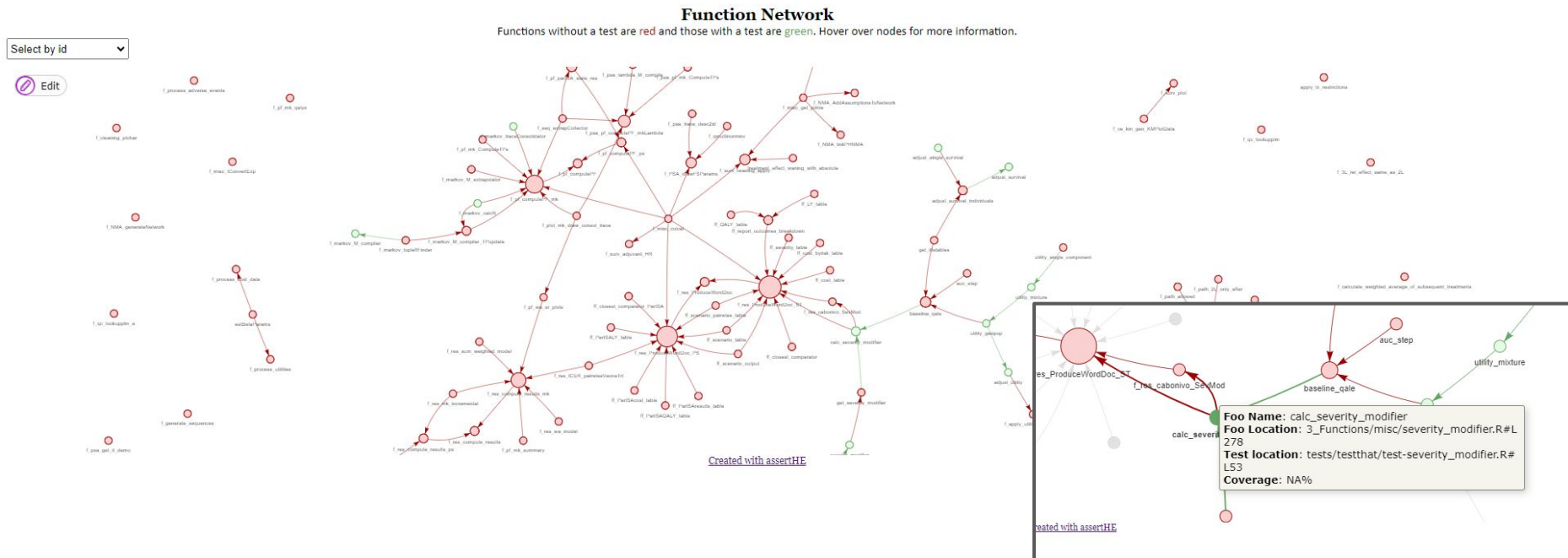
Edit



Created with assertHE

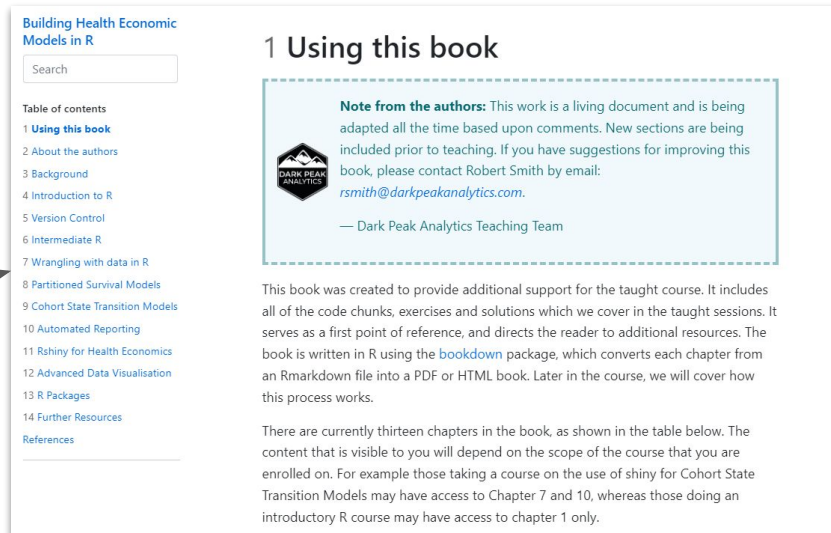
<https://github.com/DanPollardSheff/Embedding-Economic-Analysis>

Case Study: NICE RCC Pilot



Next steps

- Welcome contributions
 - Testing the software on your models
 - Suggesting improvements (see contribution page on GitHub)
 - Expansion of testing functionality
- Future development:
 - LLM Chatbot integration (using DPA teaching material to fine-tune).
 - Language selection (in progress)
- Open access publication imminent
- CRAN submission (Autumn '24)
- **Continued open-access development to maintain a collaborative tool**



The screenshot shows the book's interface. On the left is a 'Table of contents' sidebar with a search bar and a list of chapters. Chapter 1, 'Using this book', is highlighted. An arrow points from the 'LLM Chatbot integration' bullet point in the 'Next steps' list to this chapter. The main content area shows the title '1 Using this book' and a 'Note from the authors' box. The note states that the work is a living document and provides contact information for Robert Smith. Below the note, a paragraph explains the book's purpose and its creation using the bookdown package. At the bottom, another paragraph mentions that there are currently thirteen chapters and that access to certain chapters depends on the course enrollment.

Building Health Economic Models in R

Search

Table of contents

- 1 Using this book
- 2 About the authors
- 3 Background
- 4 Introduction to R
- 5 Version Control
- 6 Intermediate R
- 7 Wrangling with data in R
- 8 Partitioned Survival Models
- 9 Cohort State Transition Models
- 10 Automated Reporting
- 11 Rahiny for Health Economics
- 12 Advanced Data Visualisation
- 13 R Packages
- 14 Further Resources

References

1 Using this book

Note from the authors: This work is a living document and is being adapted all the time based upon comments. New sections are being included prior to teaching. If you have suggestions for improving this book, please contact Robert Smith by email: rsmith@darkpeakanalytics.com.

— Dark Peak Analytics Teaching Team

This book was created to provide additional support for the taught course. It includes all of the code chunks, exercises and solutions which we cover in the taught sessions. It serves as a first point of reference, and directs the reader to additional resources. The book is written in R using the [bookdown](#) package, which converts each chapter from an Rmarkdown file into a PDF or HTML book. Later in the course, we will cover how this process works.

There are currently thirteen chapters in the book, as shown in the table below. The content that is visible to you will depend on the scope of the course that you are enrolled on. For example those taking a course on the use of shiny for Cohort State Transition Models may have access to Chapter 7 and 10, whereas those doing an introductory R course may have access to chapter 1 only.

Pre-print available



https://www.courses.darkpeakanalytics.com/products/digital_downloads/asserthe-quality-assurance

My experience

1. Much, much, better at summarising (dumbing down) than problem solving.
2. Will make things up, with confidence.
3. Volume matters - very good at small tasks. Hard to articulate nuance of bigger tasks.
4. Carefully written prompts are crucial, garbage in, garbage out.
5. Replacement for human Googling + Stack Overflow.

Potential

1. More efficient model builds.
2. Freeing up time for methods development.
3. Catching bugs & errors. If review > build then particularly useful for ERGs.
4. Making models more transparent to stakeholders?
5. Cross-pollination of ideas with other disciplines.

Challenges

1. Danger for open-source (does 'no attribution' destroy thought leadership?).
2. Will make things up, with confidence.
3. If simpler tasks (e.g. writing/reviewing specific functions) are done by genAI, how do people learn & develop?
4. How do we keep the human in the loop?! (Do we need to?)



Challenges

1. **Danger for open-source (does 'no attribution' destroy thought leadership?).**
2. Will make things up, with confidence.
3. If simpler tasks (e.g. writing/reviewing specific functions) are done by genAI, how do people learn & develop?
4. How do we keep the human in the loop?! (Do we need to?)



Challenges

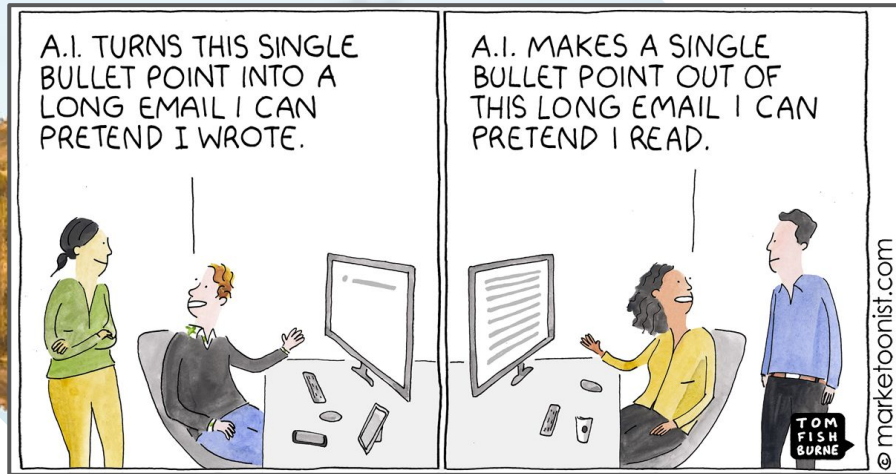
1. Danger for open-source (does 'no attribution' destroy thought leadership?).
2. **Will make things up, with confidence.**
3. If simpler tasks (e.g. writing/reviewing specific functions) are done by genAI, how do people learn & develop?
4. How do we keep the human in the loop?! (Do we need to?)



Challenges

1. Danger for open-source (does 'no attribution' destroy thought leadership?).
2. Will make things up, with confidence.
3. **If simpler tasks (e.g. writing/reviewing specific functions) are done by genAI, how do people learn & develop?**
4. How do we keep the human in the loop?! (Do we need to?)

Challenges



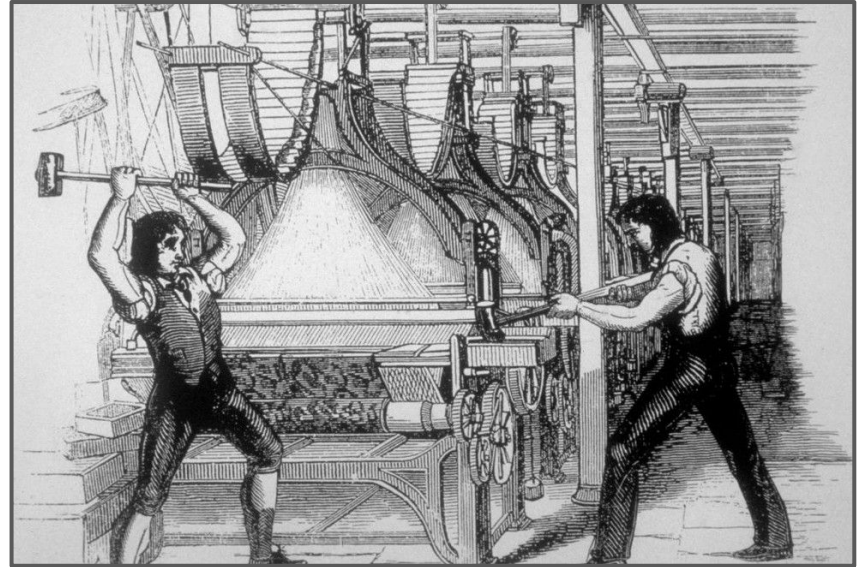
1. Danger for open-source (does 'no attribution' destroy thought leadership?).
2. Will make things up, with confidence.
3. If simpler tasks (e.g. writing/reviewing specific functions) are done by genAI, how do people learn & develop?
4. **How do we keep the human in the loop?! (Do we need to?)**

Challenges

1. Danger for open-source (does 'no attribution' destroy thought leadership?).
2. Will make things up, with confidence.
3. If simpler tasks (e.g. writing/reviewing specific functions) are done by genAI, how do people learn & develop?
4. How do we keep the human in the loop?! (Do we need to?)
5. Will we all be out of work?



Challenges



5. Will we all be out of work?

Short Courses in R for Health Economic Evaluation



Introduction to R for Health Economic Evaluation

A Dark Peak Analytics Short Course

12 | 100% peer reviewed paper of...

Introduction to R for Health Economic Evaluation

Course

This course provides an introduction to R and RStudio for those new to the language who would like a gradual introduction before progressing...



Making Health Economic Models Shiny


A Dark Peak Analytics Short Course

12 | 100% peer reviewed paper of...

Making Health Economic Models Shiny

Course

This course teaches delegates to create interactive web applications for health economic models. Based on peer reviewed literature, it shows how...



Efficient Microsimulation Modelling in R

A Dark Peak Analytics Short Course

12 | 100% peer reviewed paper of...

Efficient Microsimulation Modelling in R

Course

This course teaches delegates to build microsimulation models in base R. It shows how to use vectorization, parallelization and C++...



Automating Health Economic Evaluation with R

A Dark Peak Analytics Short Course

12 | 100% peer reviewed paper of...

Automating Health Economic Evaluation with R

Course

This course covers the skills necessary to create automated health economic evaluation reports which always reflect the latest information. Bas...



Packaging Cost-effectiveness Models in R

A Dark Peak Analytics Short Course

12 | 100% peer reviewed paper of...

Packaging Cost-effectiveness Models in R

Course

This course covers the skills necessary to build health economic evaluation models into R packages. Following our peer reviewed paper of...



R for Health Economic Evaluation Reading List

A Dark Peak Analytics Short Course

12 | 100% peer reviewed paper of...

R for Health Economic Evaluation Reading List

Course

We often get asked: "What resources exist to learn how to build health economic evaluation models in R?". We've put together this list of open-access...

